

Please replace the paragraph at p. 1, lns. 22-29 with the following paragraph:

*a<sup>2</sup>* Currently, most systems do not deal with the problem that the sampling frequency might differ considerably between the sending and the receiving side. One possible solution proposed in, EP-0680033 A2, works on pitch periods. Adding or removing pitch periods in the speech signal achieves a different duration of a speech segment without affecting other speech characteristics other than speed. This proposed solution might be used as an indirect sample rate conversion method.

Please replace the paragraph at p. 2, lns. 1-11 with the following paragraph:

*a<sup>3</sup>* Another solution uses the beginning of talkspurts as an indication to reset the playout buffer to a specified level. The distance, in number of samples, between two consecutive talkspurts is increased if the receiving side is playing faster than the sending side and decreased if the receiving side is playing slower than the sending side. In IP-telephony solutions using the IP/UDP/RTP-protocols (Internet Protocol/User Datagram Protocol/Real Time Protocol), a marker flag in the RTP header is used to identify the beginning of a talkspurt. At the beginning of a talkspurt, the playout buffer is set to a suitable size.

Please replace the paragraph at p. 2, lns. 12-20 with the following paragraph:

*a<sup>4</sup>* The solution according to EP-0680033 A2, where pitch periods are removed or inserted, assumes a fixed conversion factor between the receiving and transmitting side. Therefore, it cannot be

*a4 cont*  
used in dynamic systems, i.e. where the sampling frequencies varies. Further, it does not solve the problem with underrun or overrun situations, but is instead focused on changing the playback rate of a speech signal stored in compressed form for playback later and at a different speed to that at which it was stored.

---

Please replace the paragraph at p. 2, lns. 21-28 with the following paragraph:

---

*a5*  
Using the method of resetting the playout buffer to a certain size causes problems if there are very long talkspurts, e.g. broadcast from one speaker to several listeners. Since the length of a talkspurt is not defined in the beginning of the talkspurt, the size to reset to might be either too small or too large. If it is too small, underrun will occur and if it is too large, unnecessary delay is introduced. Thus, the problem persists.

---

Please replace the paragraph at p. 2, lns. 29-31 with the following paragraph:

---

*a6*  
The general problem with the currently known approaches is that they are static and inflexible. Therefore, dynamic solutions are required.

---

Please replace the paragraph at p. 3, lns. 8-13 with the following paragraph:

---

*a7*  
When sampling frequencies are not controlled, underrun or overrun might occur in the playout buffer at the receiving side, which causes audible artifacts in the speech signal. To avoid said overrun or underrun there is a need for dynamically keeping the playout buffer to an average size, i.e. controlling the fullness of the playout buffer.

---

Please replace the paragraph at p. 3, Ins. 14-16 with the following paragraph:

a<sup>8</sup> One object of the present invention is thus to provide a method for reducing audio artifacts in a speech signal due to overrun or underrun in the playout buffer.

Please replace the paragraph at p. 3, Ins. 17-18 with the following paragraph:

a<sup>9</sup> Another object of the invention is to dynamically control the fullness of the playout buffer so as not to introduce extra delay.

Please replace the paragraph at p. 3, Ins. 19-29 with the following paragraph:

a<sup>10</sup> The above mentioned and other objects are achieved by means of dynamic sample rate and conversion of speech frames, i.e. converting speech frames comprising N samples to instead comprise either N+1 or N-1 samples. More specifically, the invention works on an LPC-residual of the speech frame. By adding or removing a sample in the LPC-residual, a sample rate conversion will be achieved. The LPC- residual is the output from an LPC-filter, which removes the short-term correlation from the speech signal. The LPC-filter is a linear predictive coding filter where each sample is predicted as a linear combination of previous samples.

Please replace the paragraph at p.3, Ins. 30-33 through p. 4, Ins. 1-4 with the following paragraph:

a<sup>11</sup> By using the proposed sample rate conversion method, the playout buffer, of e.g. an IP-telephony terminal, can be continuously controlled with only small audio artifacts. Since the method works

a11  
cont on a sample-by-sample basis, the playout buffer can be kept to a minimum and hence no extra delay is introduced. The solution also has very low complexity, especially when the LPC-residual already is available, as is the case in e.g. a speech decoder.

---

Please replace the paragraph at p. 4, lns. 10-13 with the following paragraph:

---

a12 Although aspects of the invention have been summarised above, the method and apparatus according to the appended claims define the scope of the invention.

---

Please replace the paragraph at p. 5, lns. 5-20 with the following paragraph:

---

a13 Referring to FIG. 1, a method for improving speech quality in a communication system includes a first terminal unit TRX1 transmitting speech signals having a first sample frequency  $F_1$  and a second terminal unit TRX2 receiving said speech signals, buffering them in a playout buffer 100 with said first frequency  $F_1$  and playing out from said playout buffer with a second frequency  $F_2$ . When the buffering frequency  $F_1$  is larger than the playout frequency  $F_2$ , the playout buffer 100 will eventually be filled with samples and subsequent samples will have to be discarded. When the buffering frequency  $F_1$  is lower than the playout frequency  $F_2$ , the playout buffer will run into starvation, i.e. it will no longer have any samples to play on the output. These two problems are called overrun and underrun, respectively, and cause audible artifacts like popping and clicking sounds in the speech signal.

---

Please replace the paragraph at p. 5, lns. 21-24 with the following paragraph:

a<sup>14</sup> The above and other problems with underrun and overrun are solved by using dynamic sample rate conversion based on modifying the LPC-residual of the speech signal and will be further described with reference to FIGS. 2-8.

Please replace the paragraph at p. 6, lns. 6-14 with the following paragraph:

a<sup>15</sup> By feeding a speech frame through the LPC-filter,  $H(z)$ , the LPC-residual is found. The LPC-residual, shown in FIG. 3, contains pitch pulses  $P$  generated by the vocal cords. The distance  $L$  between two pitch pulses  $P$  is called lag. The pitch pulses  $P$  are also predictable, and since they represent the long-term correlation of the speech signal they are predicted through an LTP-filter given by the distance  $L$  between the pitch pulses  $P$  and the gain  $b$  of a pitch pulse  $P$ . The LTP-filter is usually denoted:

Please replace the paragraph at p. 6, lns. 16-19 with the following paragraph:

a<sup>16</sup> When the LPC-residual is fed through the inverse of the LTP-filter  $F(z)$ , an LTP-residual is created. In the LTP-residual, the long-term correlation in the LPC-residual is removed, giving the LTP-residual a noise-like appearance.

Please replace the paragraph at p. 6, lns. 20-27 through p. 7, lns. 1-7 with the following paragraph:

a<sup>17</sup>  
The solution according to the invention modifies the LPC-residual, shown in FIG. 3, on a sample-by-sample basis. That is, an LPC-residual block comprising N samples is converted to an LPC-residual block comprising either N+1 or N-1 samples. The LPC-residual contains less information and less energy compared to the speech signal, but the pitch pulses  $P$  are still easy to locate. When modifying the LPC-residual, samples that are close to a pitch pulse  $P$  should be avoided, because these samples contain more information and thus have a large influence on the speech synthesis. The LTP-residual is not as suitable as the LPC-residual to use for modification since the pitch pulse positions  $P$  are no longer available. Thus, the LPC-residual is better suited for modification both compared to the speech signal and to the LTP-residual, since the pitch pulses  $P$  are easily located in the LPC-residual.

Please replace the paragraph at p. 7, lns. 8-9 with the following paragraph:

a<sup>18</sup>  
A sample rate conversion consists of four modules, shown in FIG. 4:

Please replace the paragraph at p. 7, lns. 12-13 with the following paragraph:

a<sup>19</sup>  
2) LPC-Residual Extraction (LRE) modules 410 that are used to obtain the LPC-residual

$r_{LPCi}$

Please replace the paragraph at p. 7, lns. 14-18 with the following paragraph:

- a20 3) Sample Rate Conversion Methods (RCM) modules 420 that find the position at which to add or remove samples and determine how to perform the insertion and deletion, i.e. converting the LPC residual block  $r_{LPC}$  comprising N samples to a modified LPC-residual block  $r'_{LPC}$  comprising N+1 or N-1 samples; and

Please replace the paragraph at p. 7, lns. 21-23 with the following:

a21 An idea behind embodiments of the invention is that it is possible to change the playout rate of the playout buffer 440 by removing or adding samples in the LPC-residual  $r_{LPC}$ .

Please replace the paragraph at p. 7, lns. 24-27 through p. 8, lns. 1-11 with the following paragraph:

a22 The SRC module 400 decides whether samples should be added or removed in the LPC residual  $r_{LPC}$ . This is done on the basis of at least one of the four following parameters: the sampling frequencies of the sending TRX1 and receiving terminal units TRx2, information about the speech signal e.g. a voice activity detector signal, status of the playout buffer, an indicator of the beginning of a talkspurt. The four parameters are designated SRC Inputs in FIG. 4. On the basis of a function of one or several of these parameters the SRC 400 decides when to insert or remove a sample in the LPC residual  $r_{LPC}$  and optionally which RCM 420 to use. Since digital processing of speech signals usually is made on a frame-by-frame basis, the decision of when to remove or

*a<sup>22</sup>*  
*cont* add samples basically is to decide within which LPC-residual  $r_{LPC}$  frame the RCM 420 is to insert or remove a sample.

Please replace the paragraph at p. 8, lns. 12-17 with the following paragraph:

*a<sup>23</sup>* There are basically three methods of obtaining the LPC-residual  $r_{LPC}$  that is needed as input to the RCM's 420. The methods depend on the implementation of the speech encoder and will be described with reference to FIGS. 5A-5F. The LRE solution also directly influences the SSM solution, which will become apparent below.

Please replace the paragraph at p. 8, lns. 19-34 through p.9, lns. 1-4 with the following paragraph:

*a<sup>24</sup>* In FIG. 5A an analysis-by-synthesis speech encoder 500 with LTP-filter 540 is shown. This is a hybrid encoder where the vocal tract is described with an LPC-filter 550 and the vocal cords is described with an LTP-filter 540, while the LTP-residual  $\hat{r}_{LPC}^{(n)}$  is waveform-compared with a set of more or less stochastic codebook vectors from a fixed codebook 530. The input signal S is divided into frames 510 with a typical length of 10-30 ms. For each frame the LPC-filter 550 is calculated through an LPC-analysis 520 and the LPC-filter 550 is included in a closed loop to find the parameters of the LTP-filter 540. The speech decoder 580 is included in the encoder and consists of the fixed codebook 530, whose output  $\hat{r}_{LPC}^{(n)}$  is connected to the LTP-filter 540, whose output  $\hat{r}_{LPC}^{(n)}$  is connected to the LPC-filter 550, which generates an estimate  $\hat{s}(n)$  of the original speech signal  $s(n)$ . Each estimated signal  $\hat{s}(n)$  is compared with the original speech



a24  
cont  
signal  $s(n)$  and a difference signal  $e(n)$  is calculated. The difference signal  $e(n)$  is then weighted by an error-weighting block 560 to calculate a perceptual weighted error measure  $e_w(n)$ . The set of parameters that gives the least perceptual weighted error measure  $e_w(n)$  is transmitted to a receiving side 570.

Please replace the paragraph at p. 9, lns. 6-12 with the following paragraph:

a25  
As can be seen in FIG. 5C, the LPC-residual  $\hat{r}_{LPC}^{(n)}$  is the output from the LTP-filter 540.

SRC/RCM modules 545 can be connected directly to the output of the LTP-filter 540 and integrated into the speech encoder. An LRE consists of the fixed codebook 530 and the long-term predictor 540 and the SSM consists of an LPC-filter 550, thus the LRE-module and the SSM-module are natural parts of the speech decoder.

Please replace the paragraph at p. 9, lns. 13-27 with the following paragraph:

a26  
If the speech encoder, on the other hand, is an analysis-by-synthesis speech encoder where the LTP-filter 540 is exchanged to an adaptive codebook 590 as shown in FIG. 5B, the LPCresidual  $LPC(n)$  is the output from the sum of the adaptive and the fixed codebooks 590 and 530. All other elements have the same function as in FIG. 5A which shows an analysis-by-synthesis speech encoder with LTP-filter 500. As can be seen in FIG. 5D the LPC residual  $\hat{r}_{LPC}^{(n)}$  is the sum of the output from the adaptive and fixed codebook 590 and 530. The SRC/RCM modules 545 can thus again be connected directly to that output and integrated into the speech encoder as shown in

a<sup>26</sup>  
cont

FIG. 5D. The LRE consists of the adaptive and the fixed codebook 590 and 530 and the SSM consists of an LPC-filter 550, thus the LRE module and the SSM module are again natural parts of the speech decoder.

---

Please replace the paragraph at p. 9, lns. 28-33 through p. 10, lns. 1-4 with the following paragraph:

---

a<sup>27</sup>

When the speech encoder has some sort of backward adaptation, it is not feasible to make alterations in the LPC-residual since this would affect the adaptation process in a detrimental way. In FIG. 5E is shown how in these cases the parameters  $\hat{s}(n)$  from the LPC-filter 550 can be fed to an inverse LPC-filter 525 placed after the speech decoder. After the sample rate conversion has been made in the SRC/RCM modules 545 an LPC-filtering 550 is performed to reproduce the speech signal. The LRE module consists of the inverse LPC-filter 525 and the SSM module consists of the LPC-filter 550.

---

Please replace the paragraph at p. 10, lns. 5-15 with the following paragraph:

---

a<sup>28</sup>

In FIG. 5F it is shown how it is possible to produce an LPC residual  $\hat{r}_{LPC}^{(n)}$  through a full LPC analysis. The output  $\hat{s}(n)$  from the speech decoder is fed to both an LPC analysis block 520 and an LPC-inverse filter 525. After the sample rate conversion has been made in the SRC/RCM modules 545, an LPC filtering 550 is performed to reproduce the speech signal. The LRE consists in this case of the LPC analysis 520 respective the LPC inverse filter 525 and the SSM

a<sup>28</sup>  
cont module consists of the LPC filter 550. Performing an LPC analysis is considered to be well known to a person skilled in the art and is therefore not discussed any further.

---

Please replace the paragraph at p. 10, lns. 16-23 with the following paragraph:

---

a<sup>29</sup> Referring again to FIG. 4, assume that the SRC-module 400 has decided that a sample should be added or removed in the LPC residual  $r_{LPC}$  and that the LRE module 410 has produced an LPC residual  $r_{LPC}$ . The RCM-module 420 then has to find the exact position in the LPC-residual  $r_{LPC}$  where to add or remove a sample and performing the adding respective removing. There are four different methods for the RCM-module 420 to find the insertion or deletion point.

---

Please replace the paragraph at p. 10, lns. 24-28 with the following paragraph:

---

a<sup>30</sup> The first and most primitive method arbitrarily removes or adds a sample whenever this becomes necessary. If the sample rate difference between the terminals is small this will only lead to minor artifacts since the adding or removing is performed very seldom.

---

Please replace the paragraph at p. 11, lns. 17-26 with the following paragraph:

---

a<sup>31</sup> The fourth method, illustrated in FIG. 6, uses knowledge about the position  $P$  of a pitch pulse, and the lag  $L$  between two pitch pulses. With this knowledge, it is possible to calculate a position  $P'$  having low energy at which it is therefore appropriate to add or remove a sample. The new position  $P'$  can be expressed as  $P' = P + k \cdot L$ , wherein the constant  $k$  is selected so that  $P'$  is

a<sup>31</sup>  
cont selected to be somewhere in the middle between two pitch pulses, thus avoiding positions with high energy. A typical value of  $k$  is in the range of 0.5 to 0.8.

---

Please replace the paragraph at p. 11, lns. 27-31 with the following paragraph:

---

a<sup>32</sup> When the RCM-module 420 has calculated the position at which to add or remove a sample it must be determined how to perform the insertion or deletion. There are three methods of performing such insertion or deletion depending on the type of LRE-module used.

---

Please replace the paragraph at p. 12, lns. 1-7 with the following paragraph:

---

a<sup>33</sup> In the first method, either zeros are added or samples with small amplitudes are removed. This method can be used for all LRE solutions described above. (See FIGS. 5C-5F.) Notice that in FIGS. 5C and 5D the SRC/RCM-modules are placed before the synthesis filter SSM, but after the feed back of the LPC residual to the LTP-filter 540 respective the adaptive codebook 590.

---

Please replace the paragraph at p. 12, lns. 8-15 with the following paragraph:

---

a<sup>34</sup> In the second method, insertion is carried out by adding zeros and interpolating surrounding samples. Deletion is performed by removing samples and preferably smoothing surrounding samples. This method can also be used for all of the LRE solutions described above. (See FIGS. 5C-5F.) Notice that in FIGS. 5C and 5D the SRC/RCM-modules are placed before the synthesis filter SSM, but after the feed back of the LPC residual to the LTP-filter 540 respective the adaptive codebook 590.

---

Please replace the paragraph at p. 12, lns. 16-25 with the following paragraph:

In the third method, the SRC/RCM-modules 545 are placed within the feedback loop of the speech decoder instead of after the feedback loop as in the previous methods. (See FIGS. 5G-5J.)

a<sup>35</sup> Placing the SRC/RCM-modules within the feedback loop uses real LPC residual samples for the sample rate conversion, by changing the number of components in the LPC-residual. The implementation differs depending on whether it is an analysis-by-synthesis speech encoder with LTP filter shown in FIG. 5A or an analysis-by-synthesis speech encoder with adaptive codebook shown in FIG. 5B that is used.

Please replace the paragraph at p. 12, lns. 26-33 through p.13, lns. 1-2 with the following paragraph:

a<sup>36</sup> For the speech decoder with LTP filter (see FIG 5A) the SRC/RCM-modules 545 can be placed within the feedback loop in two different ways, either within the LTP feedback loop as shown in FIG. 5G or in the output from the fixed codebook 530 as shown in FIG. 5H. For the speech decoder with adaptive codebook (see FIG. 5B) the SRC/RCM can also be placed in two different ways, i.e. either before (FIG. 5J) or after, FIG. 5I, the summation of the outputs from the adaptive and the fixed codebook.

Please replace the paragraph at p. 13, lns. 3-21 with the following paragraph:

a<sup>37</sup> The alterations on the LPC residual consists of removing or adding samples just as before, but since the SRC/RCM-modules 545 are placed within the LTP feedback loop, some modifications